# Epigenetics and DNase-Seq

BMI/CS 776

www.biostat.wisc.edu/bmi776/

Spring 2017

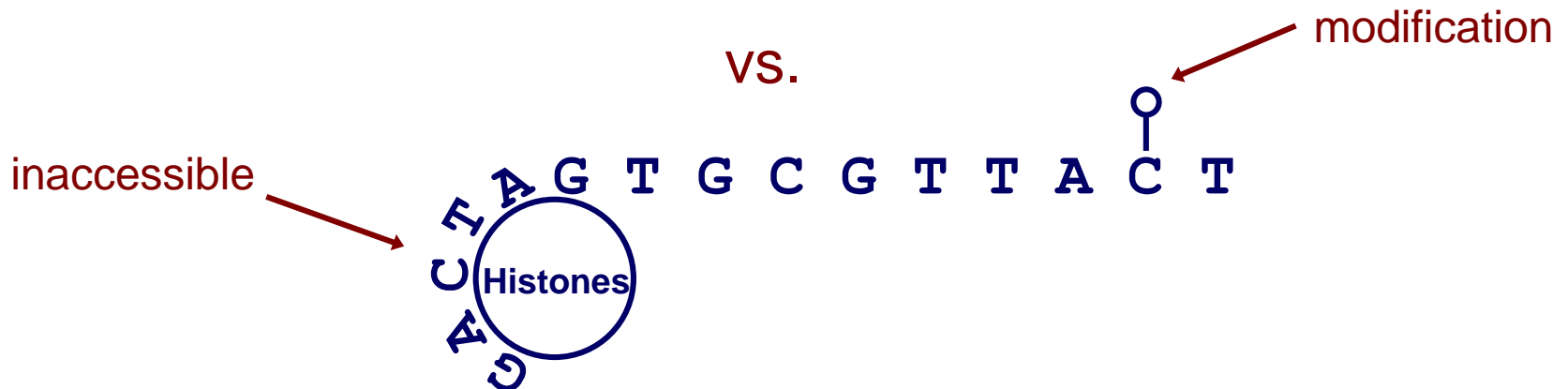Anthony Gitter

gitter@biostat.wisc.edu

# Goals for lecture

Key concepts

- Importance of epigenetic data for understanding transcriptional regulation

- Predicting transcription factor binding sites

- Gaussian Process models

# Introduction to epigenetics

# Defining epigenetics

- Formally: attributes that are "in addition to" genetic sequence or sequence modifications

- Informally: experiments that reveal the context of DNA sequence
  - DNA has multiple states and modifications

G  A  C  T  A  G  T  G  C  G  T  T  A  C  T

vs.

modification

inaccessible

**Histones**  G  T  G  C  G  T  T  A  C  T

4

# Importance of epigenetics

Better understand

- DNA binding and transcriptional regulation
- Differences between cell and tissue types
- Development and other important processes
- Non-coding genetic variants (next lecture)

# PWMs are not enough

- Genome-wide motif scanning is imprecise

- Transcription factors (TFs) bind < 5% of their motif matches

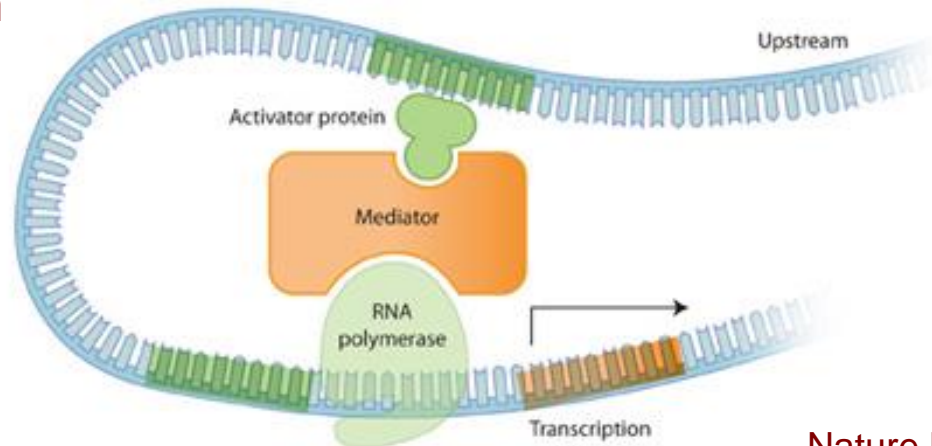- Same motif matches in all cells and conditions

# PWMs are not enough

- DNA looping can bring distant binding sites close to transcription start sites
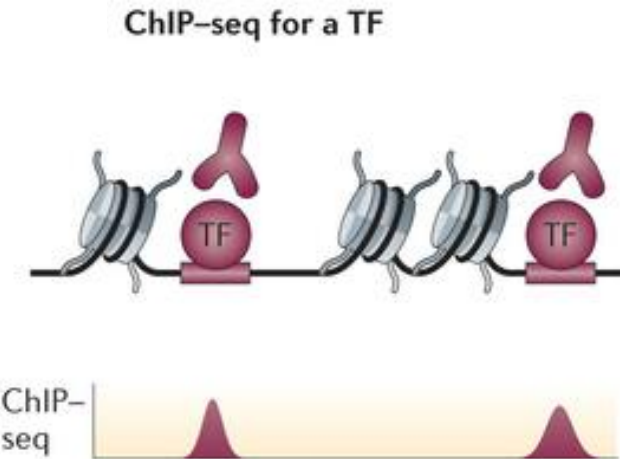- Which genes does an enhancer regulate?



Enhancer: DNA binding site for TFs, can be far from affected gene

Promoter: DNA binding site for TFs, close to gene transcription start site
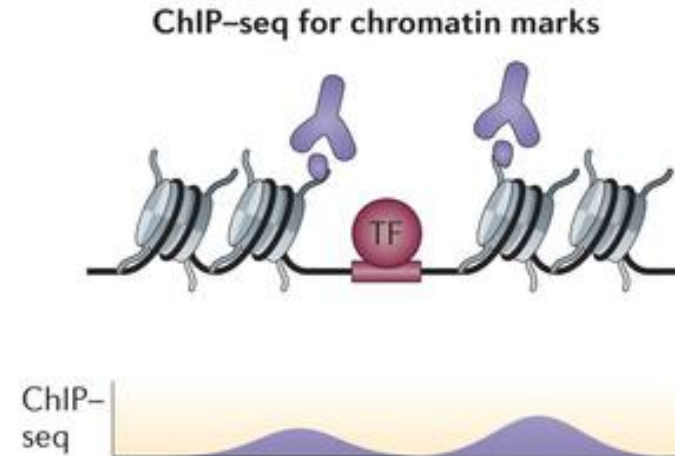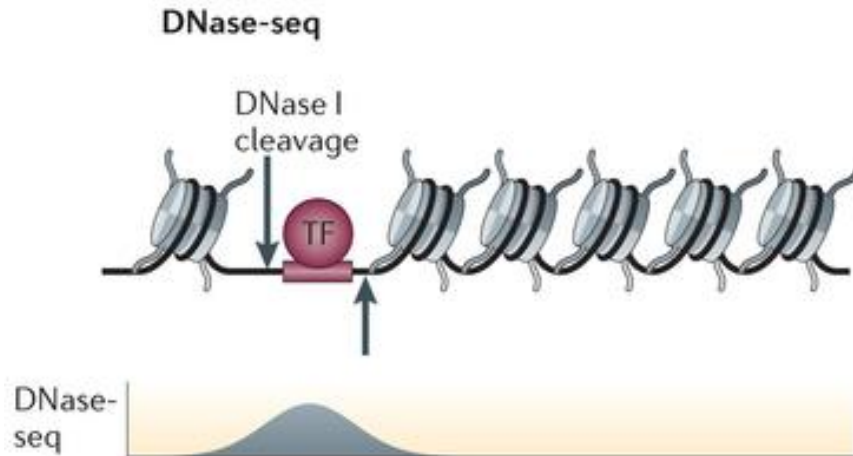
Nature Education 2010

7

# Mapping regulatory elements genome-wide

- Can do much better than motif scanning with additional data

- ChIP-seq measures binding sites for one TF at a time

- Epigenetic data suggests where *some* TF binds

ChIP–seq for a TF



ChIP–seq

Shlyueva *Nature Reviews Genetics* 2014

DNase-seq

DNase I cleavage



DNase-seq
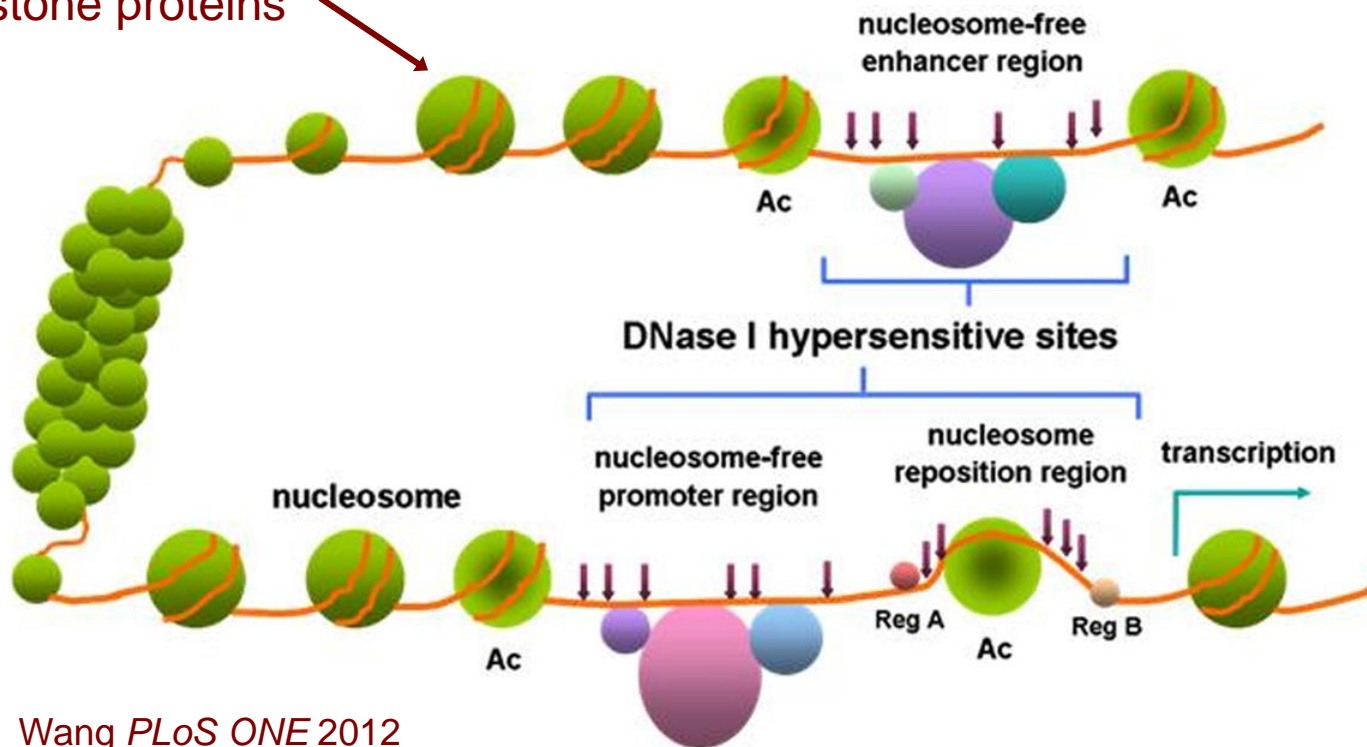
ChIP–seq for chromatin marks



ChIP–seq

# DNase I hypersensitivity

- Regulatory proteins bind accessible DNA
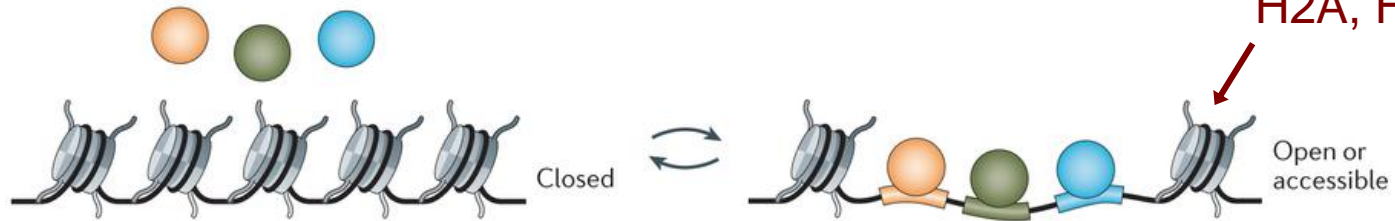- DNase I enzyme cuts open chromatin regions that are not protected by nucleosomes

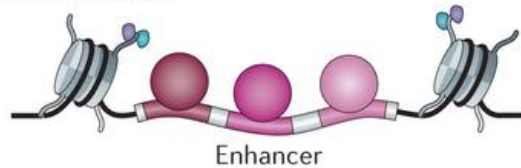Nucleosome: DNA wrapped around histone proteins



Wang *PLoS ONE* 2012
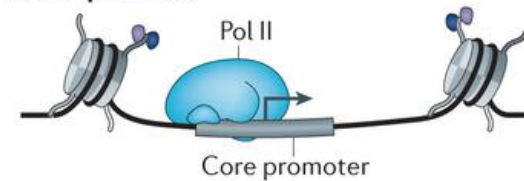
# Histone modifications

- Mark particular regulatory configurations

Two copies of histone proteins H2A, H2B, H3, H4

Chromatin as accessibility barrier

Closed

Open or accessible

Active enhancer

Enhancer

Active promoter

Pol II

Core promoter

Shlyueva *Nature Reviews Genetics* 2014

TFs     DNA binding motifs     DNA-binding proteins: TFs, CTCF, repressors and polymerases     H3K4me1     H3K27ac     H3K4me3     H3K27me3

- H3 (protein) K27 (amino acid) ac (modification)

H3   N–A R T K Q T A R K S T G G K A P R K Q L A T K A A R K S A...G V K K...–C
           2   4       8 9         14      17 18         23      26 27         36

Latham *Nature Structural & Molecular Biology* 2007; Katie Ris-Vicari

10

# DNA methylation
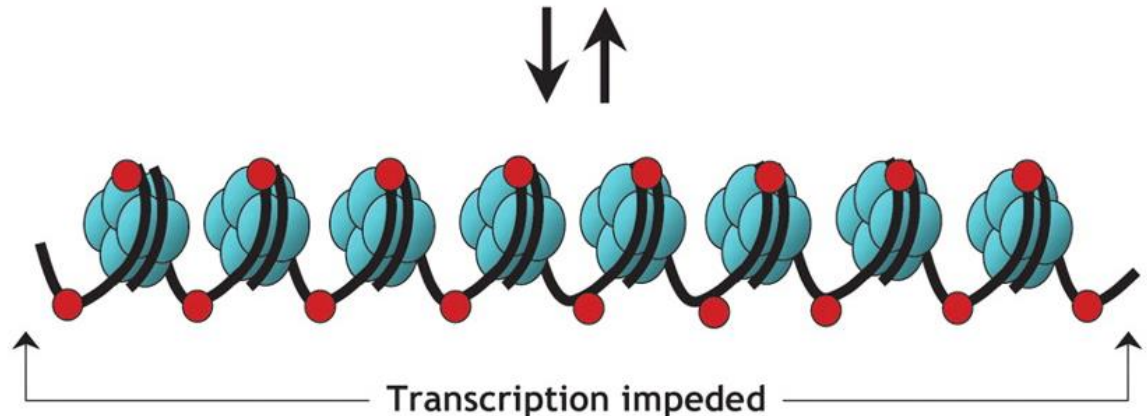
- Reversible DNA modification
- Represses gene expression

Nucleosome

DNA methylation

**Transcription possible**

Gene "switched on"
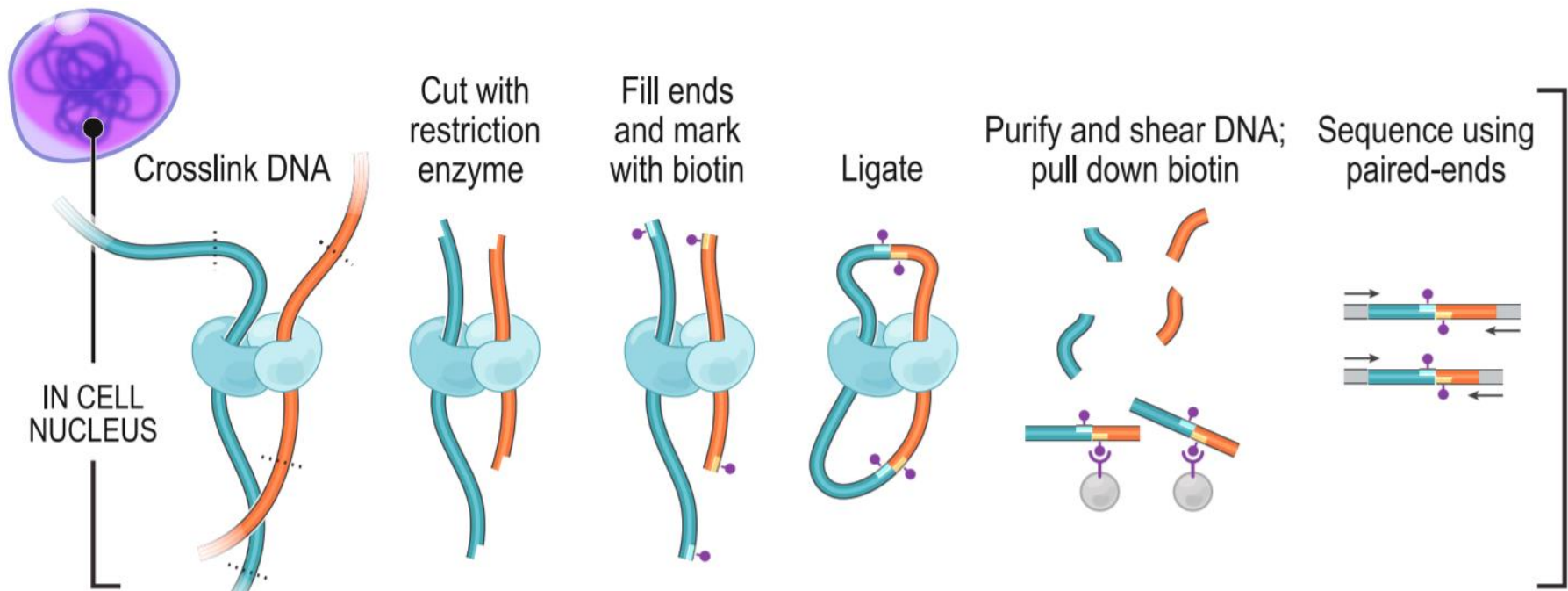- Active (open) chromatin
- Unmethylated cytosines (white circles)
- Acetylated histones

Gene "switched off"
- Silent (condensed) chromatin
- Methylated cytosines (red circles)
- Deacetylated histones
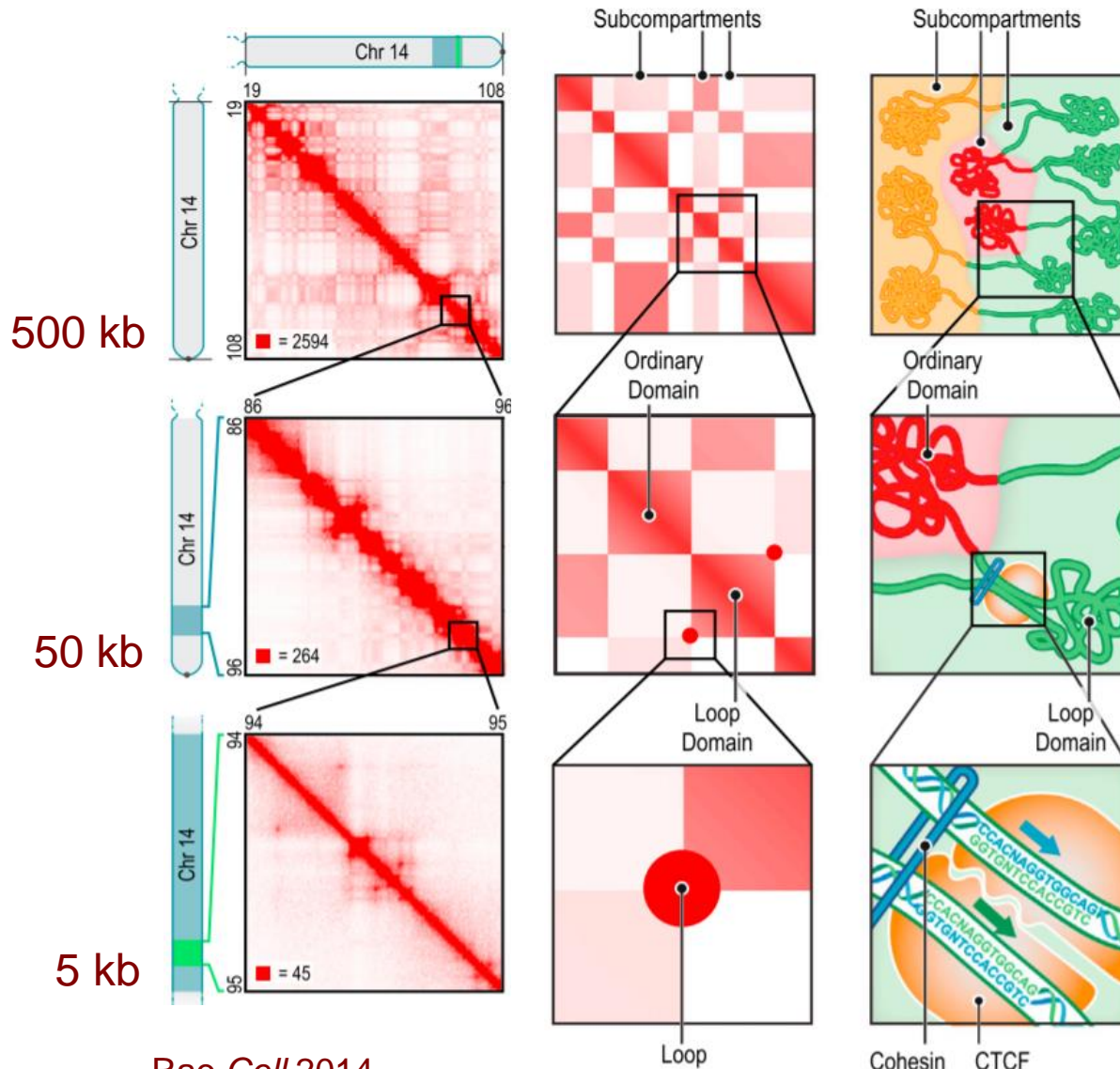
**Transcription impeded**

11

# 3d organization of chromatin

- Algorithms to predict long range enhancer-promoter interactions

- Or measure with chromosome conformation capture (3C, Hi-C, etc.)
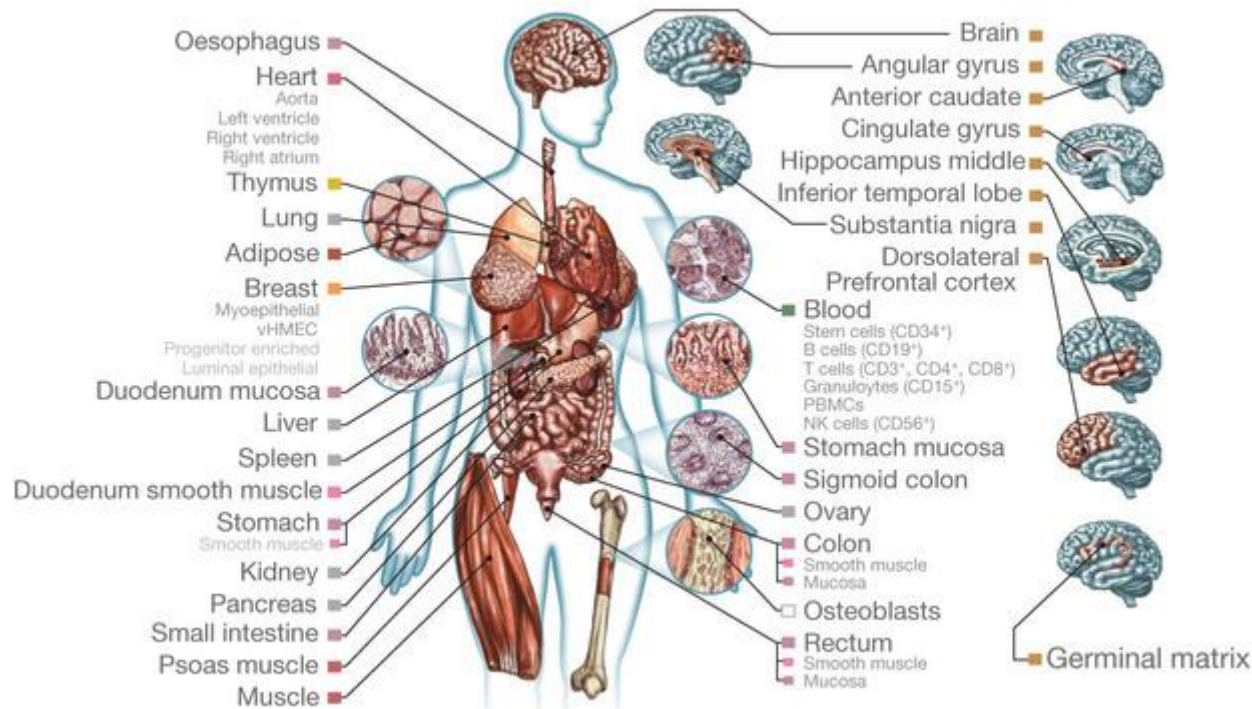


Rao *Cell* 2014

# 3d organization of chromatin



500 kb

50 kb

5 kb

Rao *Cell* 2014

- Hi-C produces 2d chromatin contact maps

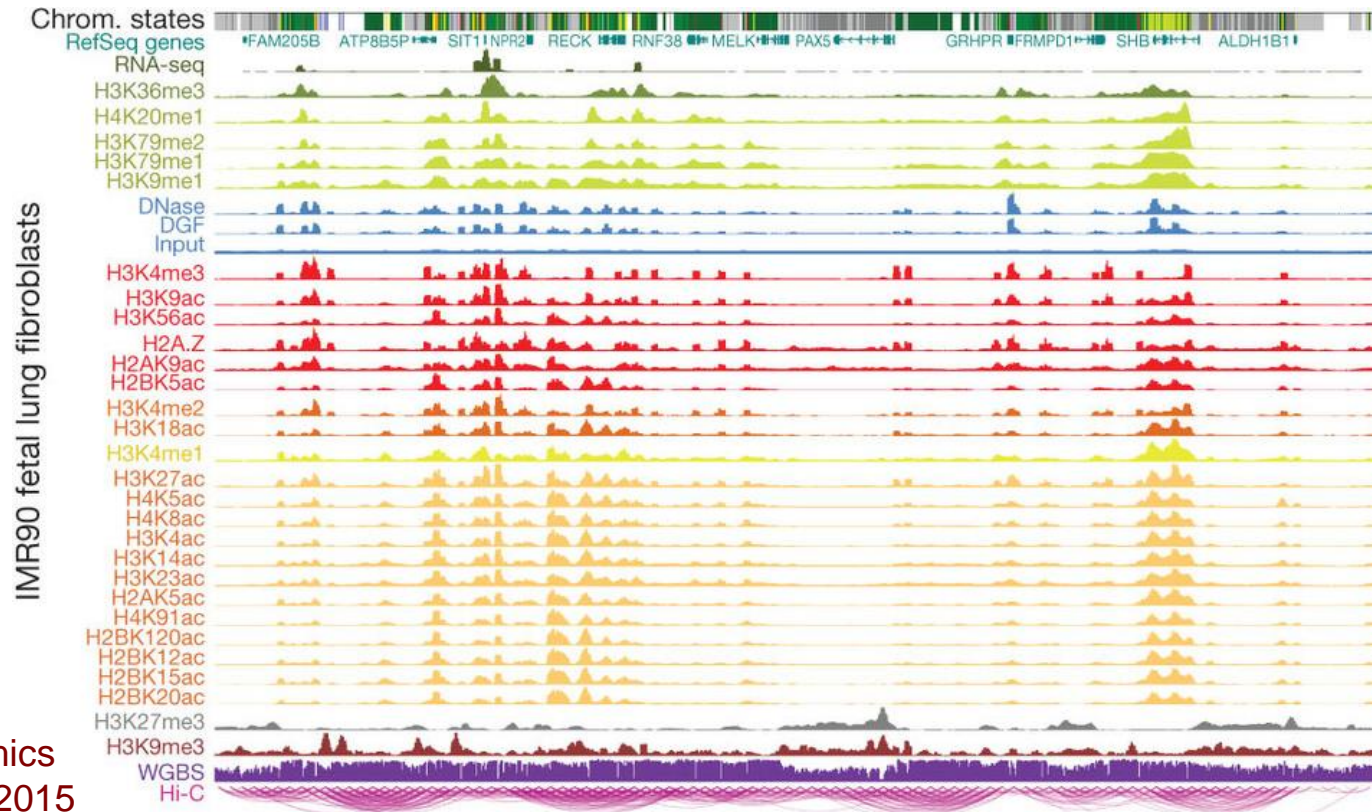- Learn domains, enhancer-promoter interactions

13

# Large-scale epigenetic maps

- Epigenomes are condition-specific
- Roadmap Epigenomics Consortium and ENCODE surveyed over 100 types of cells and tissues
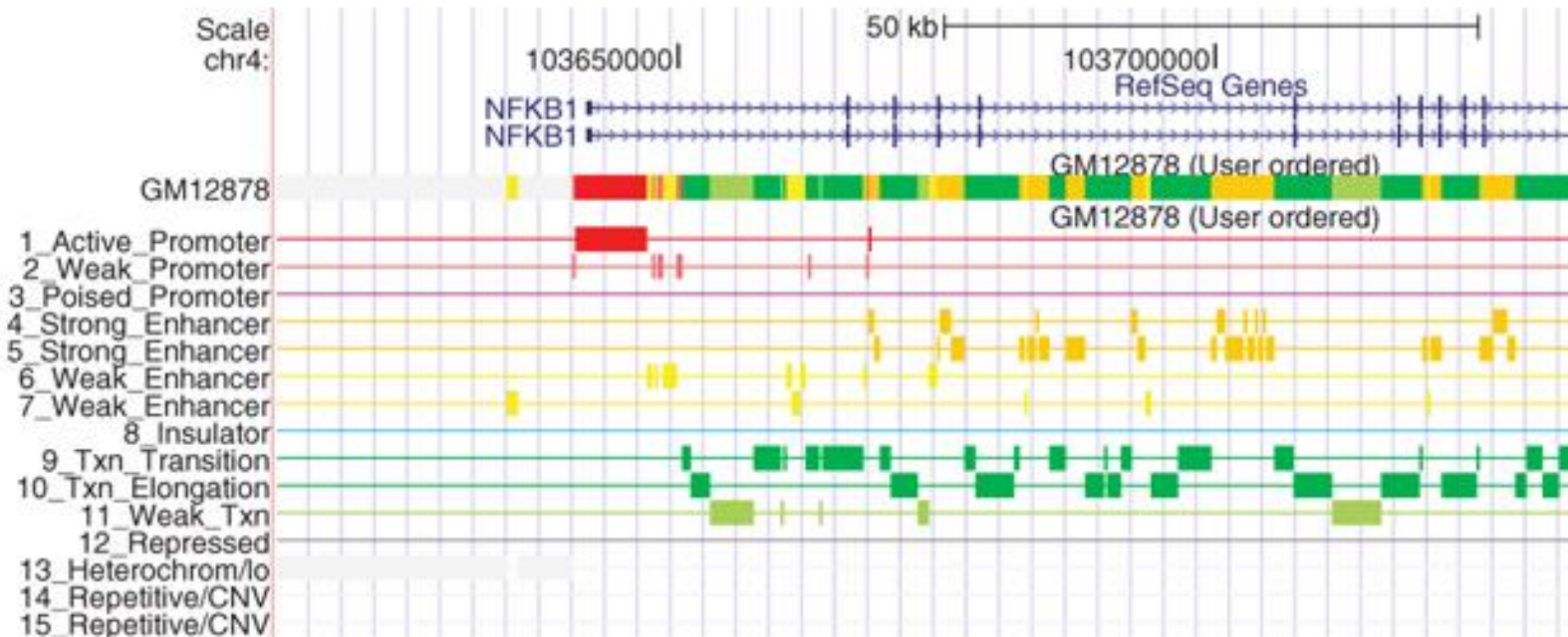
# Genome annotation

- ## Combinations of epigenetic signals can predict functional state

  - ### ChromHMM: Hidden Markov model

  - ### Segway: Dynamic Bayesian network

# Genome annotation
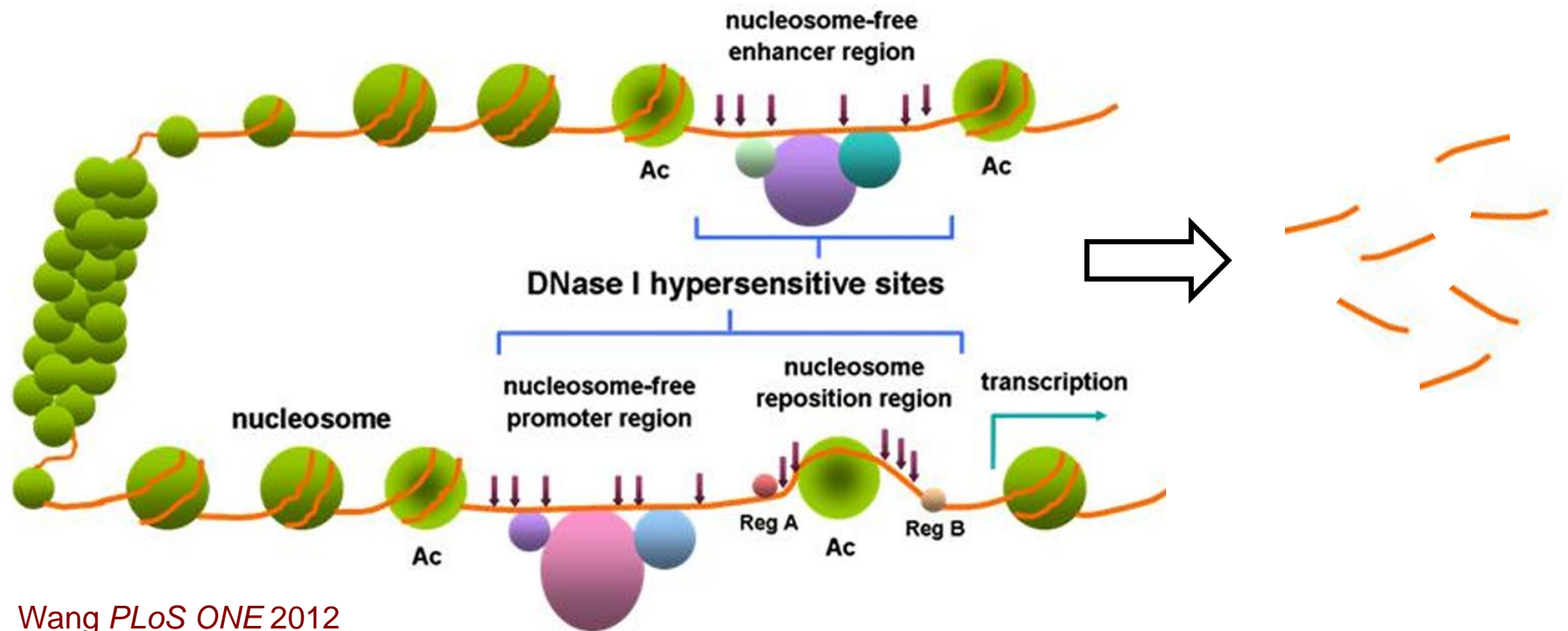
- States are more interpretable than raw data



Ernst and Kellis *Nature Methods* 2012
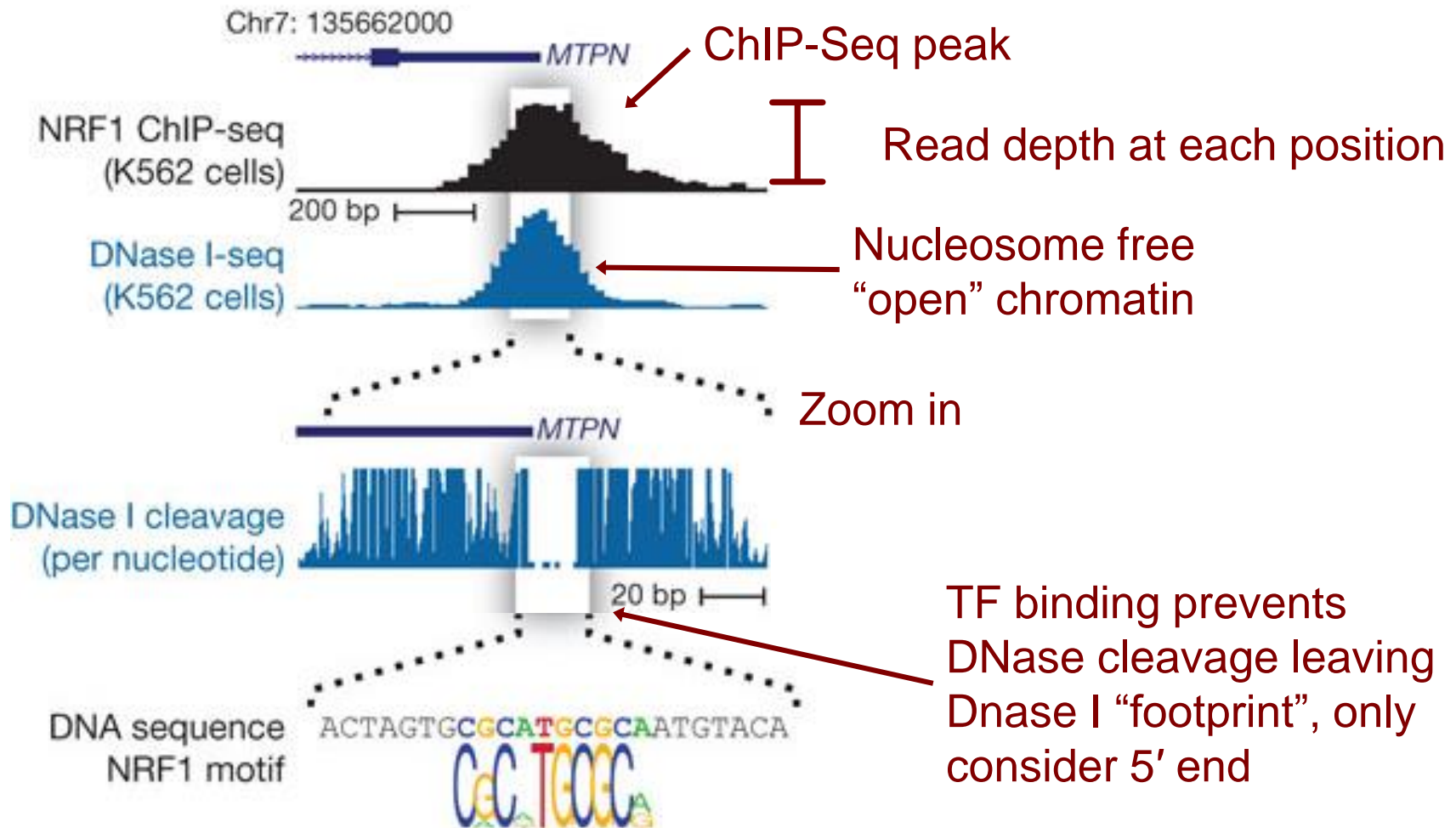
# Predicting TF binding with DNase-Seq

# DNase I hypersensitive sites

- Arrows indicate DNase I cleavage sites
- Obtain short reads that we map to the genome



Wang *PLoS ONE* 2012

# DNase I footprints

- Distribution of mapped reads is informative of open chromatin and specific TF binding sites



Chr7: 135662000

MTPN

ChIP-Seq peak

NRF1 ChIP-seq
(K562 cells)

200 bp

Read depth at each position

DNase I-seq
(K562 cells)

Nucleosome free
"open" chromatin

Zoom in

MTPN

DNase I cleavage
(per nucleotide)

20 bp

TF binding prevents
DNase cleavage leaving
Dnase I "footprint", only
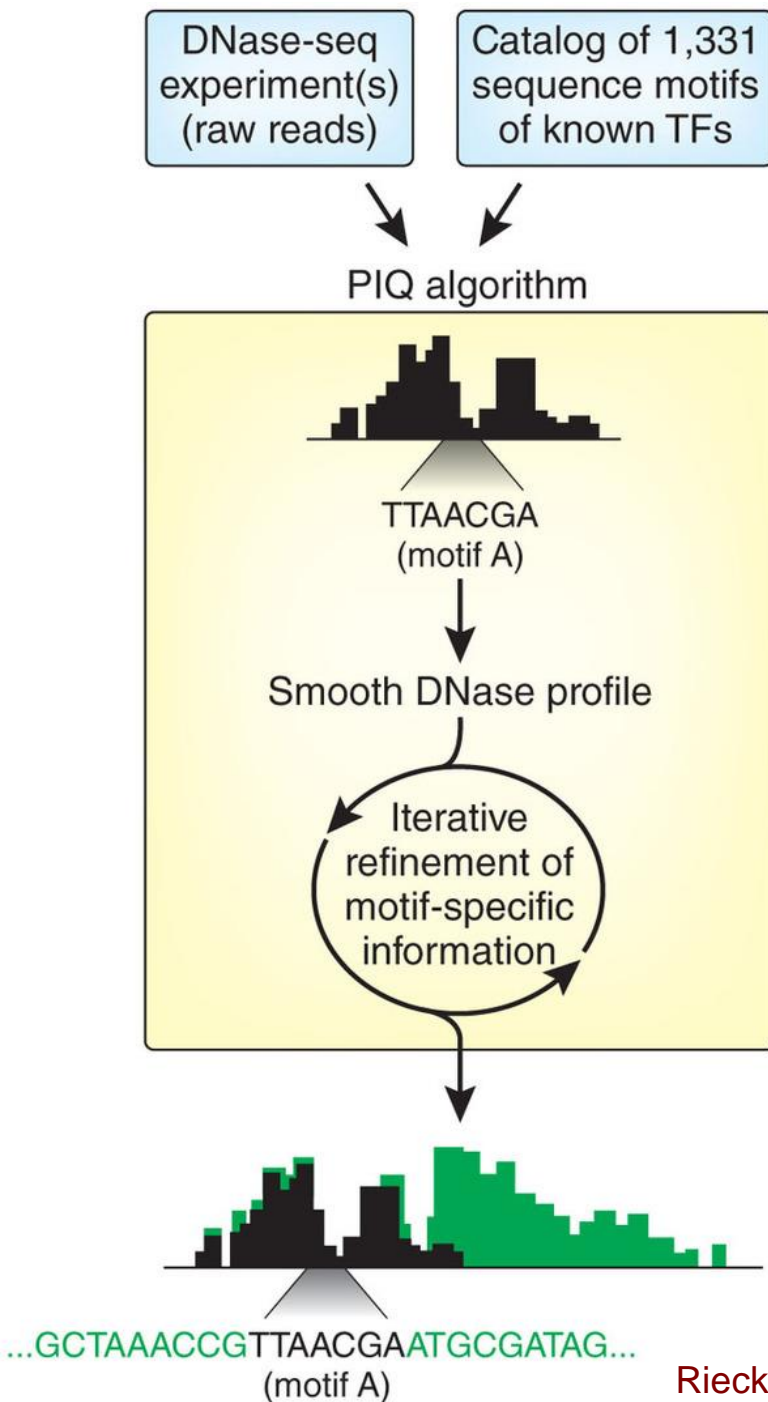consider 5′ end

DNA sequence
NRF1 motif

ACTAGTGCGCATGCGCAATGTACA

# DNase I footprints to TF binding predictions

- DNase footprints suggest that *some* TF binds that location

- We want to know *which* TF binds that location

- Two ideas:
  – Search for DNase footprint patterns, then match TF motifs
  – Search for motif matches in genome, then model proximal DNase-Seq reads

We'll consider this approach
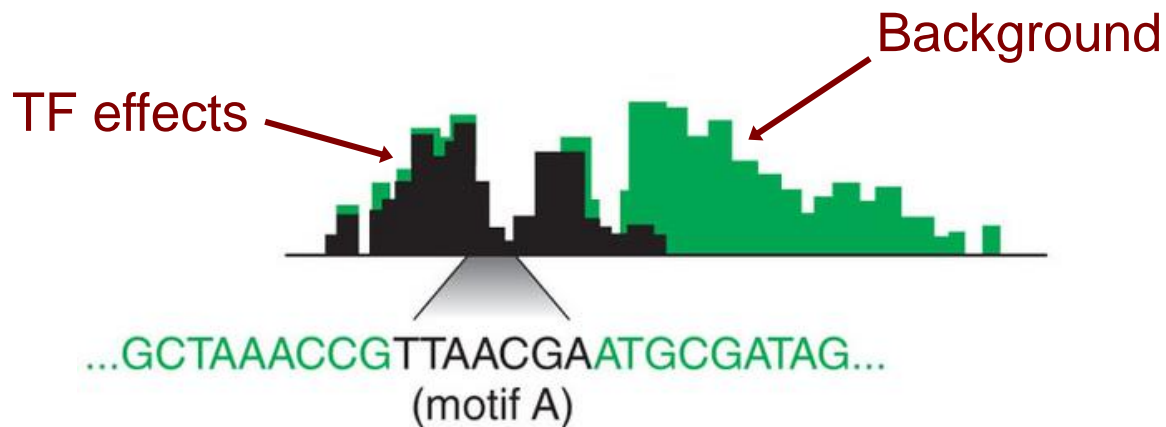
# Protein Interaction Quantification (PIQ)



...GCTAAACCGTTAACGAATGCGATAG...
(motif A)

- Sherwood et al. *Nature Biotechnology* 2014

- **Given**: TF motifs and DNase-Seq reads

- **Do**: Predict binding sites of each TF

# PIQ main idea

- With no TF binding, DNase-Seq reads come from some background distribution

- TF binding changes read density in a *TF-specific* way
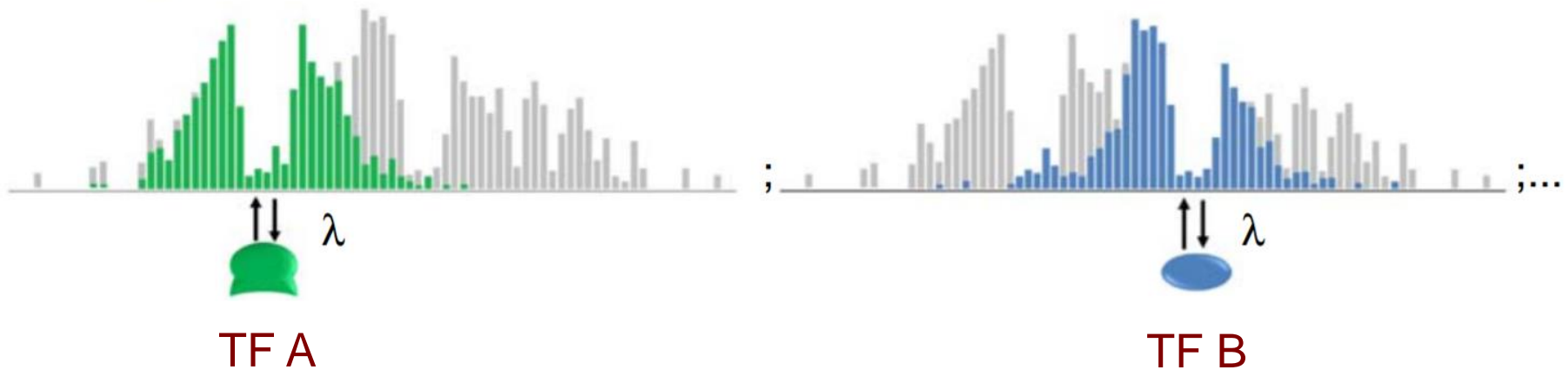
Background

TF effects

...GCTAAACCGTTAACGAATGCGATAG...
(motif A)

# PIQ main idea

- Shape of DNase peak and footprint depend on the TF



TF binding estimation

TF A          TF B

Sherwood *Nature Biotechnology* 2014

# PIQ features

- ## We'll discuss
  - Modeling the DNase-Seq background distribution
  - How TF binding impacts that distribution
  - Priors on TF binding

- ## We'll skip
  - Modeling multiple replicates or conditions, cross-experiment and cross-strand effects
  - Expectation propagation
  - TF hierarchy: pioneers, settlers, migrants

# Algorithm preview

- Identify candidate binding sites with PWMs
- Build a probabilistic model of the DNase-Seq reads
- Estimate TF binding effects
- Estimate which candidate binding sites are bound
- Predict pioneer, settler, and migrant TFs

# DNase-Seq background

- Each replicate is noisy, don't want to over-interpret this noise
  - Only counting density of 5′ ends of reads


- Manage two competing objectives
  - Smooth some of the noise
  - Don't destroy base pair resolution signal

# Gaussian processes

- Can model and smooth sequential data

- Bayesian approach

- [Jupyter notebook demonstration](Jupyter notebook demonstration)

# TF DNase profile

- Adjust the log-read rate by a TF-specific effect at binding sites

DNase profile for factor *l*

Whether site *m* is bound

$$\widehat{\mu}_l = \mu_i + \begin{cases} \beta_{i-j,l} & |y_m - j| \leq W \ and \ I_m = 1 \\ 0 & otherwise \end{cases}$$

DNase log-read rate adjusted for binding of factor *l*

Location of binding site *m*

Window size

DNase log-read rate at position *i* from Gaussian process

# TF DNase profile

- DNase profiles represented as a vector for each TF

DNase profile for factor *l*

$$\widehat{\mu_l} = \mu_i + \begin{cases} \beta_{i-j,l} & |y_m - j| \leq W \ and \ I_m = 1 \\ 0 & otherwise \end{cases}$$

Can't be too far apart

$y_m$     $i$



$\mu$    $\beta$
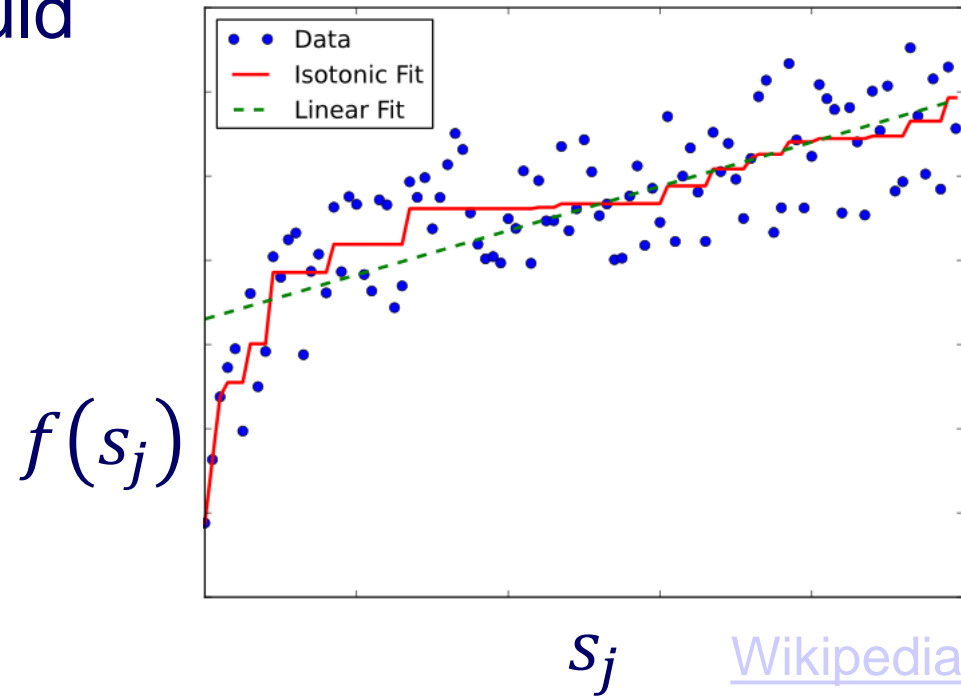
... ...

ACTAGTGCGCATGCGCAATGTACA

$W$     $l =$     $W$

# Priors on TF binding

- TF binding event $I_j$ should be more likely when
  - motif score $s_j$ is high
  - DNase counts $c_j$ are high

- Isotonic (monotonic) regression

Legend:
- Data
- Isotonic Fit
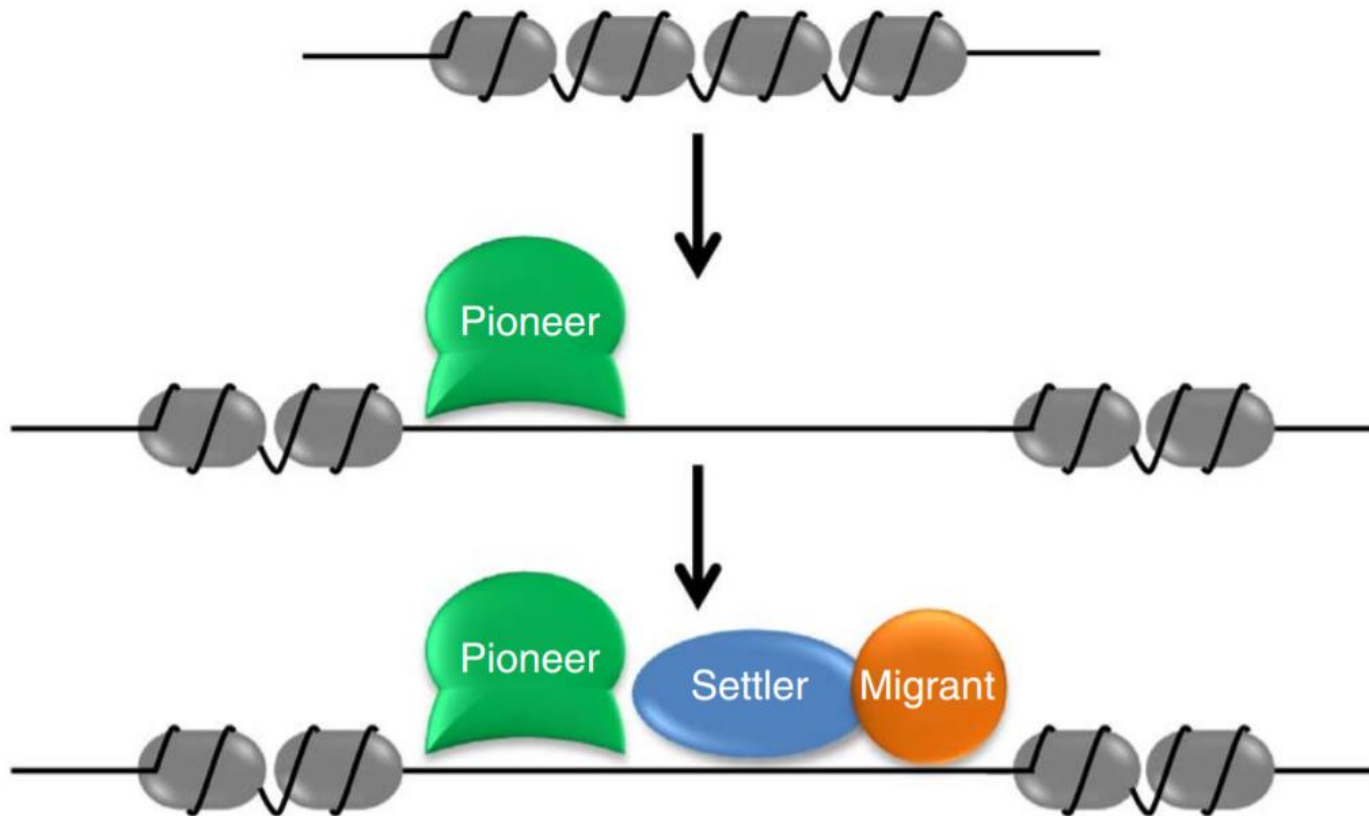- Linear Fit

$f(s_j)$

$s_j$    Wikipedia

$$\log(P(I_j = 1)) = f(s_j) + g(c_j)$$

# Full algorithm

- **Given**: TF motifs and DNase-Seq reads
- **Do**: Predict binding sites of each TF

- Identify candidate binding sites with PWMs
- Fit Gaussian process parameters for background
- Estimate TF binding effects $\beta_{i-j,l}$
- Iterate until parameters converge
  - Estimate Gaussian process posterior with expectation propagation
  - Estimate expectation of which candidate binding sites are bound
  - Update monotonic regression functions for binding priors

# TF binding hierarchy

- Pioneer, settler, and migrant TFs
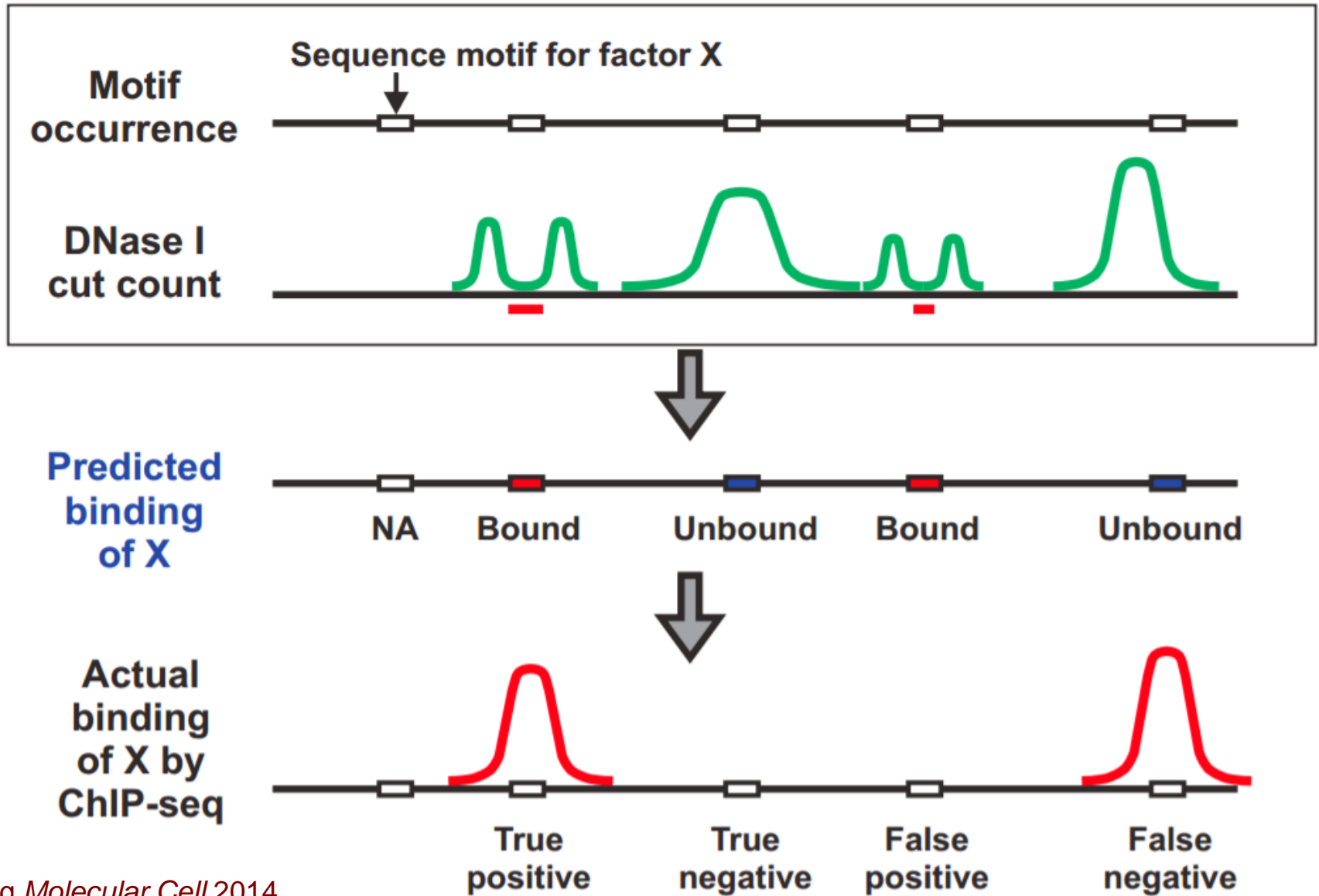


Sherwood *Nature Biotechnology* 2014

# Evaluation: confusion matrix

- Compare predictions to actual ground truth (gold standard)



Lever *Nature Methods* 2016

# Evaluation: ChIP-Seq gold standard

# Evaluation: ROC curve

- Calculate **r**eceiver **o**perating **c**haracteristic curve (ROC)
- True Positive Rate versus False Positive Rate
- Summarize with **a**rea **u**nder ROC **c**urve (AUC ROC)
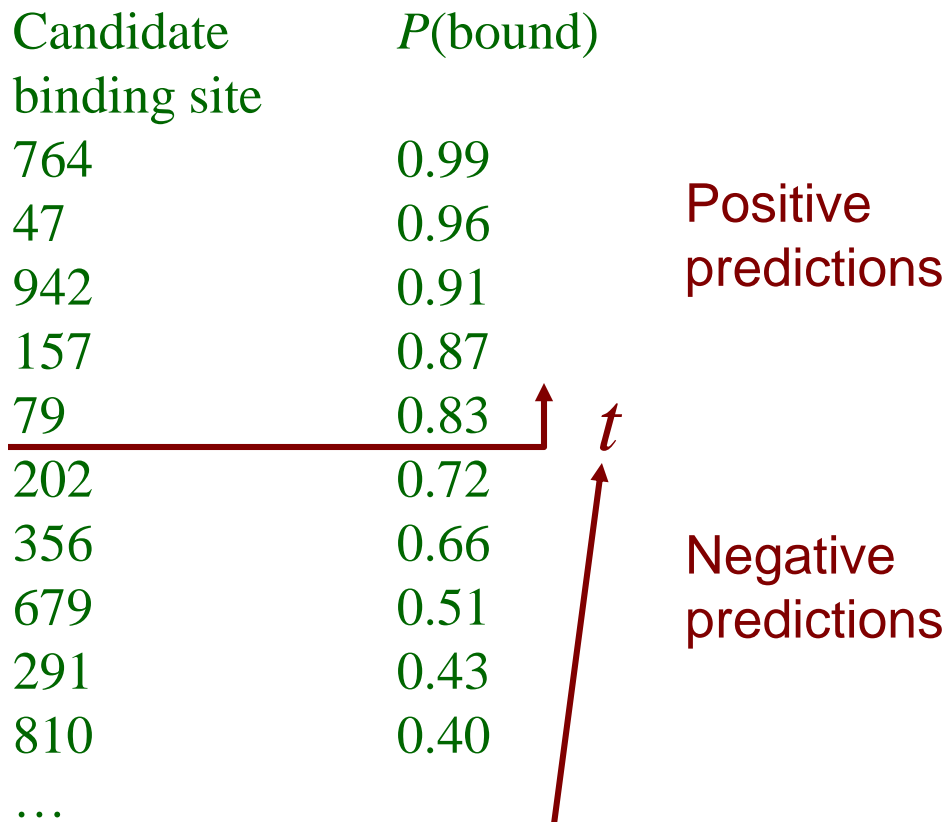
$$TPR = \frac{TP}{P} = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{N} = \frac{FP}{FP + TN}$$

Includes true negatives
Reason to prefer precision-recall for class imbalanced data
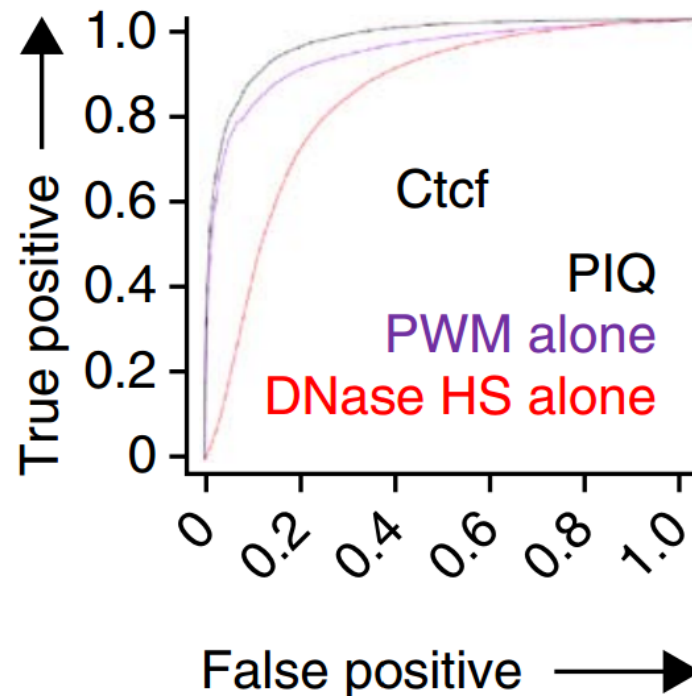
# Evaluation: ROC curve

- TPR and FPR are defined for a **set** of positive predictions

- Need to threshold continuous predictions

- Rank predictions

- ROC curve assesses all thresholds

| Candidate binding site | $P$(bound) |
|---|---|
| 764 | 0.99 |
| 47 | 0.96 |
| 942 | 0.91 |
| 157 | 0.87 |
| 79 | 0.83 |
| 202 | 0.72 |
| 356 | 0.66 |
| 679 | 0.51 |
| 291 | 0.43 |
| 810 | 0.40 |
| … | |

$t$

Positive predictions

Negative predictions
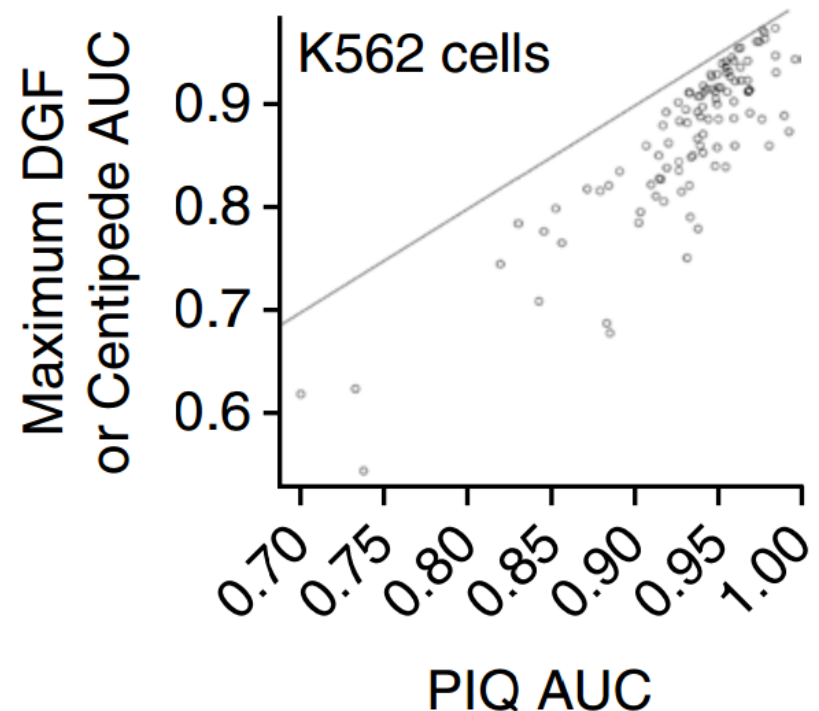
Calculate TPR and FPR at all thresholds $t$

# PIQ ROC curve for mouse Ctcf

- Compare predictions to ChIP-Seq
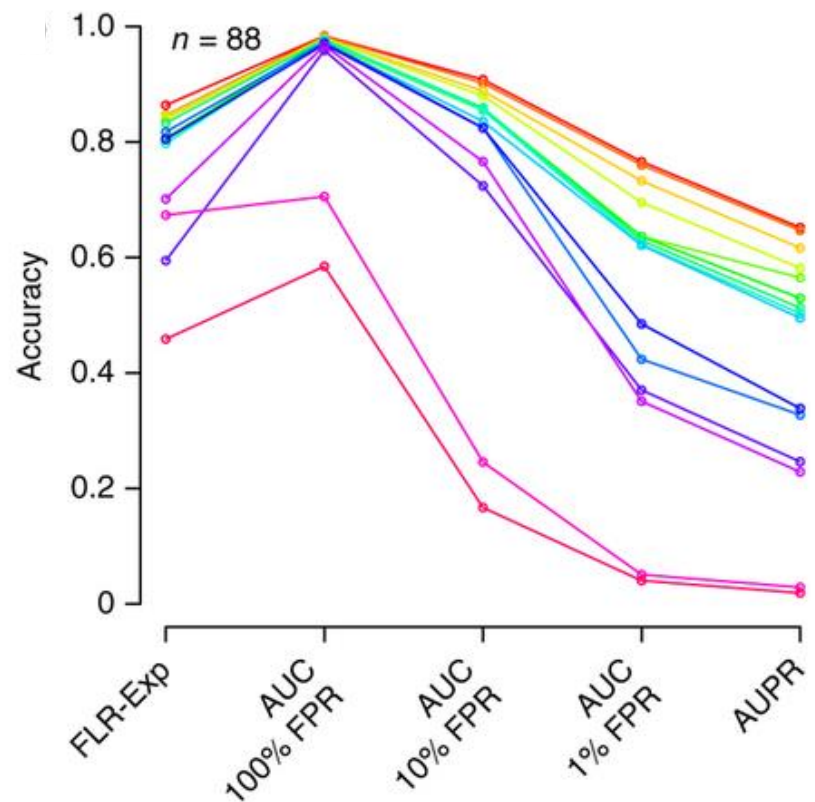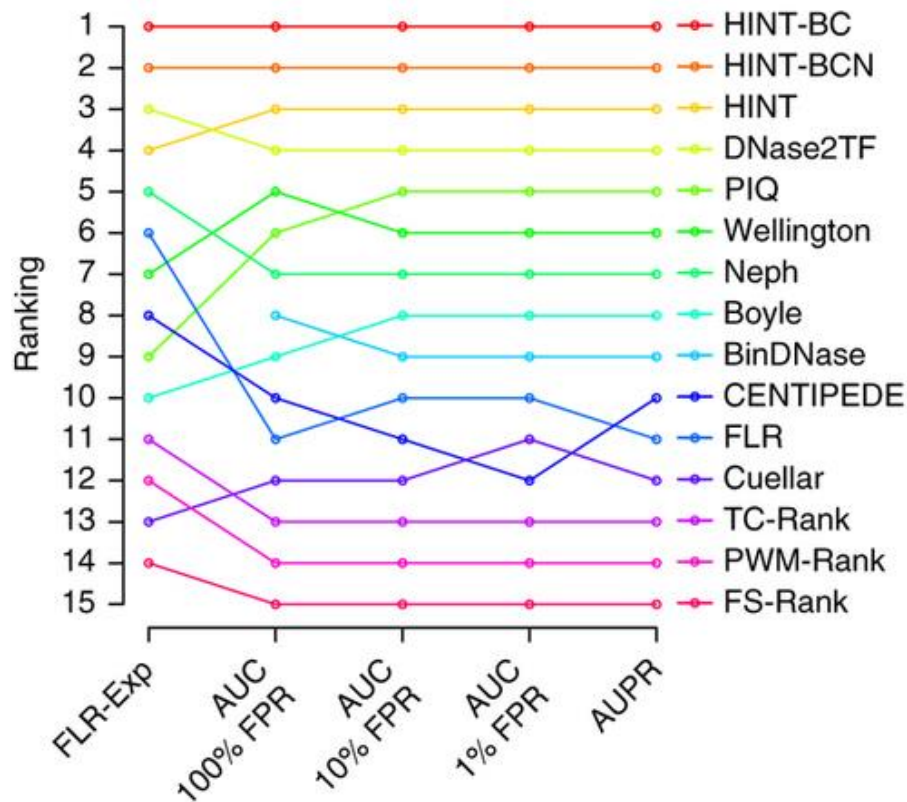- Full PIQ model improves upon motifs or DNase alone

# PIQ evaluation

- Compare to two standard methods
  - 303 ChIP-Seq experiments in K562 cells
  - Centipede, digital genomic footprinting

- Compare AUC ROC
  - PIQ has very high AUC
  - Mean 0.93
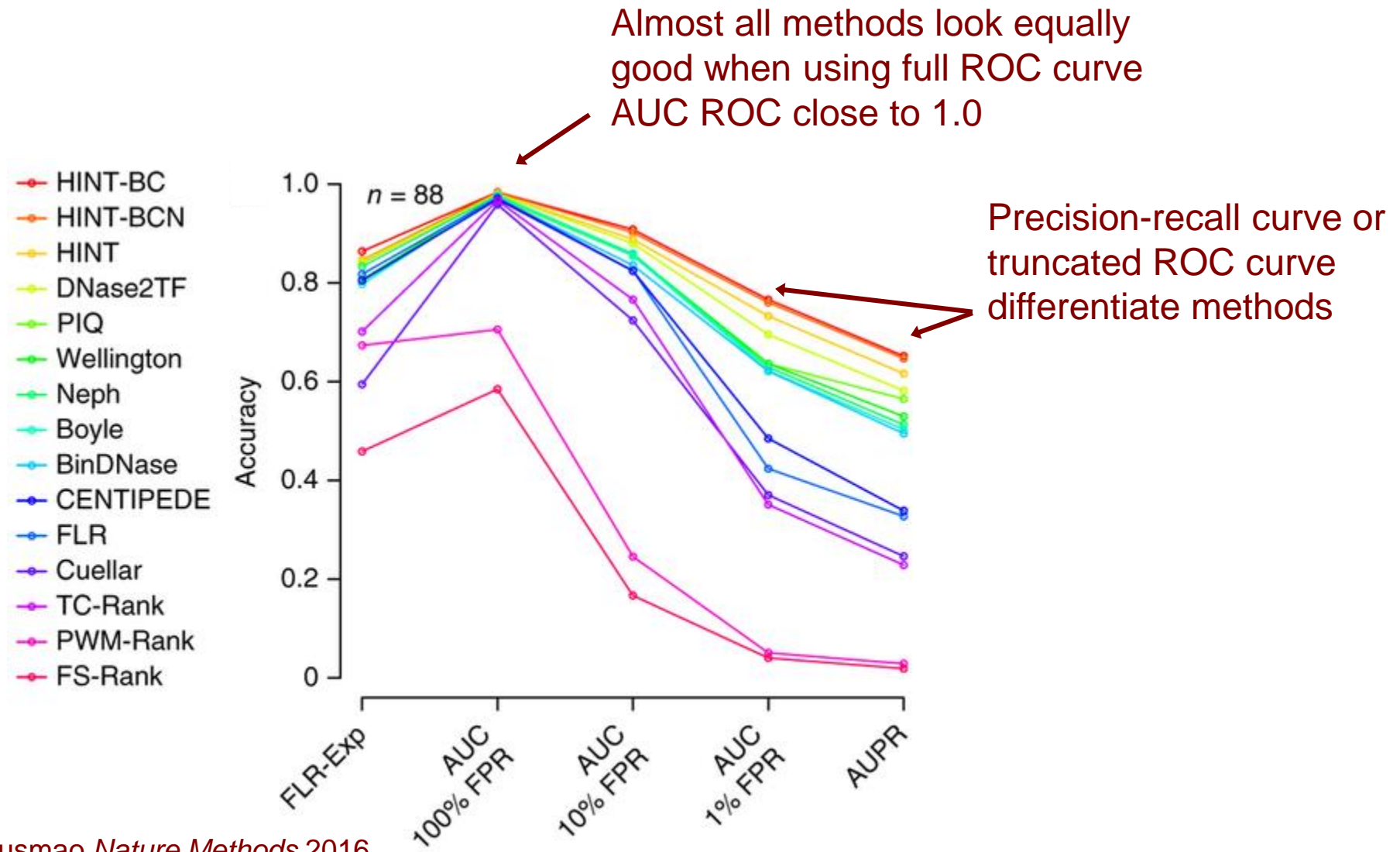  - Corresponds to recovering median of 50% of binding sites

# DNase-Seq benchmarking

- PIQ among top methods in large scale DNase benchmarking study
- HMM-based model HINT was top performer

# Downside of AUC ROC for genome-wide evaluations

Almost all methods look equally good when using full ROC curve AUC ROC close to 1.0

Precision-recall curve or truncated ROC curve differentiate methods

# PIQ summary

- Smooth noisy DNase-Seq data without imposing too much structure

- Combine DNase-Seq and motifs to predict condition-specific binding sites

- Supports replicates and multiple related conditions (e.g. time series)